

# wCorr Formulas

*Paul Bailey, Ahmad Emad, (people who do QC for this product)*

*2016-03-09*

This document shows the methodology used in computing the correlations in wCorr for the weighted Pearson, Spearman, polyserial, and polychoric correlations. For the polyserial and polychoric correlations the coefficient is estimated using a numerical likelihood maximization. The weighted likelihood functions that are used are motivated and described.

Additionally, for the polyserial and polychoric evidence, is offered for the correctness of the methods, including a sketch of a proof of the consistency of both methods and numerical simulations to show:

- The bias of the methods as a function of the true correlation coefficient ( $\rho$ ) and the number of observations ( $n$ )
- The root mean square error (RMSE) of the methods as a function of  $\rho$  and  $n$  in the weighted and unweighted cases

The *wCorr Arguments* vignette describes the effects the **ML** and **fast** arguments have on computation.

## Methodology

Here we focus on measurement of the correlation coefficients between two vectors of random variables that are jointly bivariate normal—call the vectors  $\mathbf{X}$  and  $\mathbf{Y}$ .<sup>1</sup> The  $i^{th}$  members of the vectors are then called  $x_i$  and  $y_i$ .

### Pearson and Spearman methodology

The Pearson correlation is computed using the formula

$$\rho_{Pearson} = \frac{\sum_i [w_i(x_i - \bar{x}) \times (y_i - \bar{y})]}{\sqrt{\sum_i (w_i \times (x_i - \bar{x})^2) + \sum_i (w_i \times (y_i - \bar{y})^2)}}$$

where  $\bar{x} = \frac{1}{\sum w_i} \sum_i w_i x_i$ ,  $\bar{y} = \frac{1}{\sum w_i} \sum_i w_i y_i$ , and  $n$  is the number of elements in  $\mathbf{X}$  and  $\mathbf{Y}$ .<sup>2</sup>

The Spearman correlation coefficient is calculated by first taking the rank of the data before the same formula is applied. When data are ranked ties must be handled in some way. The chosen method is to assign the average of all tied ranks. For example, if the second and third rank units are tied then both units would receive a rank of 2.5 (the average of 2 and 3).

### Polyserial methodology

For the polyserial correlation, it is again assumed that there are two continuous variables  $\mathbf{X}$  and  $\mathbf{Y}$  that have a bivariate normal distribution.<sup>3</sup>

---

<sup>1</sup>The Spearman correlation coefficient can be motivated by a much more broad class of random variables but that is not dealt with in this document.

<sup>2</sup>See the “correlate” function in Stata Corp, Stata Statistical Software: Release 8. College Station, TX: Stata Corp LP, 2003.

<sup>3</sup>For a more complete treatment of the polyserial correlation, see Cox, N. R., “Estimation of the Correlation between a Continuous and a Discrete Variable” *Biometrics*, **50** (March), 171-187, 1974.

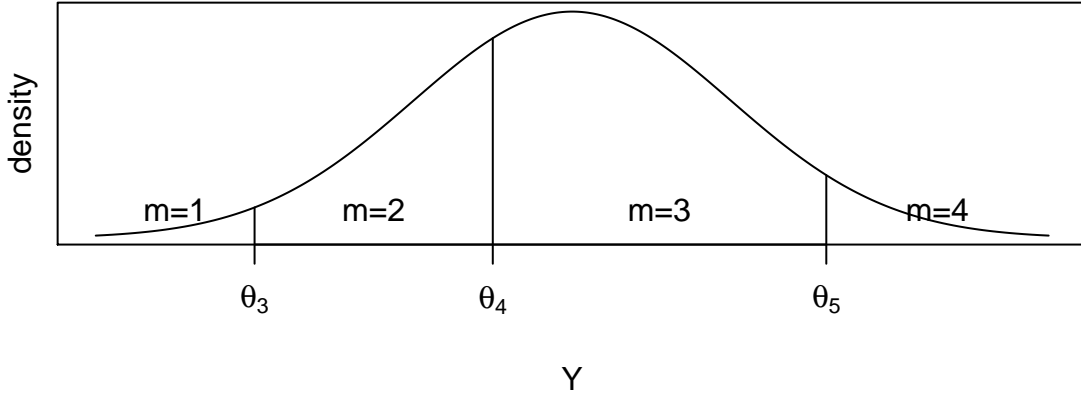
$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N \left[ \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \Sigma \right]$$

Where  $N(a, \Sigma)$  is a bivariate normal distribution with mean  $a$  and variance  $\Sigma$ . For the polyserial,  $\mathbf{Y}$  is discretized into the random variable  $\mathbf{M}$  according to

$$m_i = \begin{cases} 1 & \text{if } \theta_2 < y_i < \theta_3 \\ 2 & \text{if } \theta_3 < y_i < \theta_4 \\ \vdots & \\ t & \text{if } \theta_{t+1} < y_i < \theta_{t+2} \end{cases}$$

where, for notational convenience,  $\theta_2 \equiv -\infty$  and  $\theta_{t+2} \equiv \infty$ .<sup>4</sup>

The following figure shows the density of  $\mathbf{Y}$  when  $\theta = (-\infty, -2, -0.5, 1.6, \infty)$ . Here, for example, any value of  $-2 < y_i < -0.5$  would have  $m_i = 2$ .



Notice that  $\mu_y$  is not identified (or is irrelevant) because setting  $\tilde{\mu}_y = \mu_y + a$  and  $\tilde{\theta} = \theta + a$  lead to exactly the same values of  $\mathbf{M}$  and so one of the two must be arbitrary assigned. A convenient decision is to decide  $\mu_y \equiv 0$ . A similar argument holds for  $\sigma_y$  so that  $\sigma_y \equiv 1$ .

For  $\mathbf{X}$  Cox (1974) observes that the MLE mean and standard deviation of  $\mathbf{X}$  are simply the average and (population) standard deviation of the data and do not depend on the other parameters.<sup>5</sup> This can be taken advantage of by defining  $z \equiv \frac{x - \bar{x}}{\hat{\sigma}_x}$ .

Combining these simplifications the probability of any given  $x_i, m_i$  pair is

$$\Pr(\rho = r; Z = z_i, M = m_i) = \phi(z_i) \int_{\theta_{m_i+1}}^{\theta_{m_i+2}} dy f(y|Z = z, \rho = r)$$

where  $\Pr(\rho = r; Z = z, M = m)$  is the probability of the event  $\rho = r$  given the data  $z$  and  $m$ ,  $\phi(\cdot)$  is the standard normal and  $f(y|z, \rho)$  is the distribution of  $y$  conditional on  $z$  and  $\rho$ . Because  $y$  and  $z$  are jointly

<sup>4</sup>The indexing is somewhat odd to be consistent with Cox (1974). Nevertheless, this treatment does not use the Cox definition of  $\theta_0, \theta_1$  or  $\theta_2$  which are either not estimated (as is the case for  $\theta_0$ , and  $\theta_1$ ) or are reappropriated (as is the case for  $\theta_2$ ). Cox calls the correlation coefficient  $\theta_2$  while this document uses  $\rho$  and uses  $\theta_2$  to store  $-\infty$  as a convenience so that the vector  $\theta$  includes the (infinite) bounds as well as the interior points.

<sup>5</sup>The population standard deviation is used because it is the MLE for the standard deviation. Notice that, while the sample variance is an unbiased estimator of the variance and the population variance is not an unbiased estimator of the variance, they are very similar and the variance is also a nuisance parameter, not a parameter of interest when finding the correlation.

bivariate normally distributed (by assumption)

$$f(y|Z = z, \rho = r) = N\left(\mu_y + \frac{\sigma_y}{\sigma_z}r(z - \mu_z), (1 - r^2)\sigma_y\right)$$

because both  $\mathbf{Z}$  and  $\mathbf{Y}$  are standard normals

$$f(y|Z = z, \rho = r) = N(r \cdot z, (1 - r^2))$$

now, define  $w \equiv \frac{y - r \cdot z}{1 - r^2}$  and  $w$  has a standard normal distribution. Plugging this in

$$\Pr(\rho = r, \theta; Z = z_i, M = m_i) = \phi(z_i) \left[ \Phi\left(\frac{\theta_{m_i+2} - r \cdot z_i}{\sqrt{1 - r^2}}\right) - \Phi\left(\frac{\theta_{m_i+1} - r \cdot z_i}{\sqrt{1 - r^2}}\right) \right]$$

Where  $\Phi(\cdot)$  is the standard normal cumulative density function. Using the above probability function as an objective, the log-likelihood is then maximized.

$$\ell(\rho; z, m) = \sum_i w_i \ln [\Pr(\rho = r; Z = z_i, M = m_i)]$$

where  $w_i$  is the weight of the  $i^{th}$  case.

The value of the nuisance paramter  $\theta$  is chosen to be  $\Phi^{-1}(n/N)$  where  $n$  is the number of values to the left of the cut point ( $\theta_i$  value) and  $N$  is the number of data points overall.

## Polyserial computation

The derivatives of  $\ell$  can be computed but are not readily computed and so when the ML argumet is set to **FALSE** (the default) a one dimensional optimization of  $\rho$  is calculated using `stats::optimize`. When the ML argument is set to **TRUE** a multi-dimensional optimization is done for  $\rho$  and  $\theta$  using `minqa::bobyqa`. As is shown below, the difference between these two is slight, if present, and so the default value of ML is recommended.

Because the optimization is not perfect when the correlation is in a boundary condition ( $\rho \in \{-1, 1\}$ ), a check for perfect correlation is performed before the above optimization by simply seeing if the values of  $\mathbf{X}$  and  $\mathbf{M}$  have exactly the same order.

## Polychoric methodology

Similar to the polyserial, the polychoric is a simple case of two continuous variables  $\mathbf{X}$  and  $\mathbf{Y}$  that have a bivariate normal distribution. In the case of the polyserial the continuous (latent) variable  $\mathbf{Y}$  was observed as a discretized variable  $\mathbf{M}$ . For the polychoric this is again true but now the continuous (latent) variable  $\mathbf{X}$  is observed as a discrete variable  $\mathbf{P}$  according to

$$p_i = \begin{cases} 1 & \text{if } \theta'_2 < x_i < \theta'_3 \\ 2 & \text{if } \theta'_3 < x_i < \theta'_4 \\ \vdots & \\ t & \text{if } \theta'_{t+1} < x_i < \theta'_{t+2} \end{cases}$$

where  $\theta$  remains the cut points for the distribution defining the transformation of  $\mathbf{Y}$  to  $\mathbf{M}$  and  $\theta'$  is the cut points for the transformation from  $\mathbf{X}$  to  $\mathbf{P}$ . Similar to  $\theta$ ,  $\theta'$  has  $\theta'_2 \equiv -\infty$  and  $\theta'_{t+2} \equiv \infty$ .

Similar to in the polyserial,  $\mu_y$  is not identified (or is irrelevant) because setting  $\tilde{\mu}_y = \mu_y + a$  and  $\tilde{\theta} = \theta + a$  lead to exactly the same values of  $\mathbf{M}$  and so one of the two must be arbitrary assigned. The same is true

for  $\mu_x$ . A convenient decision is to decide  $\mu_y = \mu_x \equiv 0$ . A similar argument holds for  $\sigma_y$  and  $\sigma_x$  so that  $\sigma_y = \sigma_x \equiv 1$

Then the probability of any given  $m_i, p_i$  pair is

$$\Pr(\rho = r; P = p_i, M = m_i) = \int_{\theta'_{p_i+1}}^{\theta_{p_i+2}} dx \int_{\theta_{m_i+1}}^{\theta_{m_i+2}} dy f(x, y | \rho = r)$$

where  $\rho$  is the correlation coefficient.

Using this function as an objective, the log-likelihood is then maximized.

$$\ell(\rho; p, m) = \sum_i w_i \ln [\Pr(\rho = r; P = p_i, M = m_i)]$$

### Polychoric computation

This again mirrors the treatment of the polyserial. The derivatives of  $\ell$  can be computed but are not readily computed and so when the ML argument is set to **FALSE** (the default) a one dimensional optimization of  $\rho$  is calculated using `stats::optimize`. When the ML argument is set to **TRUE** a multi-dimensional optimization is done for  $\rho, \theta$ , and  $\theta'$  using `minqa::bobyqa`. As is shown below, the difference between these two is slight, if present, and so the default value of ML is recommended.

Because the optimization is not perfect when the correlation is in a boundary condition ( $\rho \in \{-1, 1\}$ ), a check for perfect correlation is performed before the above optimization by simply seeing if the values of  $\mathbf{P}$  and  $\mathbf{M}$  have a Goodman-Kruskal correlation coefficient of -1 or 1. When this is the case, the MLE of -1 or 1, respectively, is returned.

### Correctness

It is easy to prove the consistency of the  $\theta$  for the polyserial and  $\theta$  and  $\theta'$  using the non-ML case. Similarly, for  $\rho$ , because it is an MLE that can be obtained by taking a derivative and setting it equal to zero, the results are asymptotically unbiased and obtain the Cramer-Rao lower bound.

This does not speak to the small sample properties of these correlation coefficients. Previous work has described their properties by simulation and so that tradition is continued below.

### General setup for the unweighted case

A simulation is run several times. For each iteration, the following procedure is used:

- select the number of observations ( $n$ )
- select a true correlation coefficient  $\rho$
- generate  $\mathbf{X}$  and  $\mathbf{Y}$  to be bivariate normally distributed using a pseudo-Random Number Generator (RNG)
- using a pseudo-RNG, select the the number of bins for  $\mathbf{M}$  and  $\mathbf{P}$  ( $t$  and  $t'$ ) independantly from the set  $\{2, 3, 4, 5\}$
- select the bin boundaries for  $\mathbf{M}$  and  $\mathbf{P}$  ( $\theta$  and  $\theta'$ ) by sorting the results of  $(t - 1)$  and  $(t' - 1)$  draws, respectively, from a normal distribution using a pseudo-RNG
- confirm that at least 2 levels of each of  $\mathbf{M}$  and  $\mathbf{P}$  are occupied (if not, rerun to generating  $\mathbf{X}$  and  $\mathbf{Y}$ )
- calculate and record relevant statistics

When the exact method of selecting a parameter (such as  $n$ ) is not noted in the above description it is described as part of each simulation.

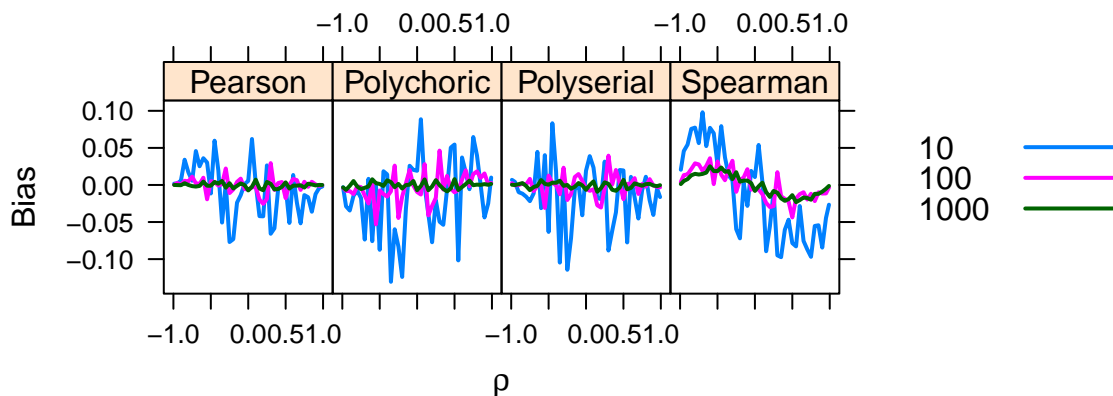
## bias of the correlations

This sections shows the bias of the correlations as a function of the true correlation coefficient,  $\rho$ . To that end, a simulation was done at each level of the cartesian product of  $\rho \in (-0.99, -0.95, -0.90, -0.85, \dots, 0.95, 0.99)$ , and  $n \in \{10, 100, 1000\}$ . For precision, each iteration is run fifty times. The bias is the mean difference between the true correlation coefficient ( $\rho_i$ ) and estimate correlation coefficient ( $r$ ). the RMSE is the squareroot of the mean square error.

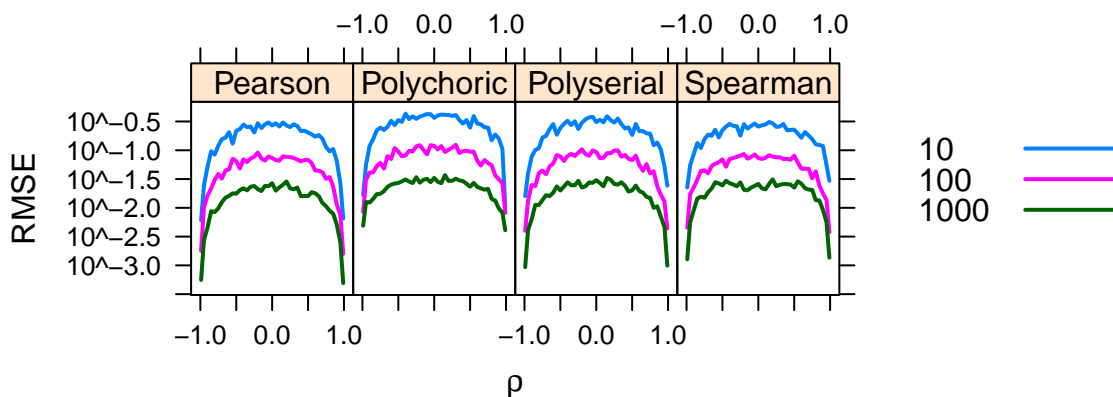
$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_i (r_i - \rho_i)^2}$$

## RMSE of the correlations as a function of $n$ and $\rho$

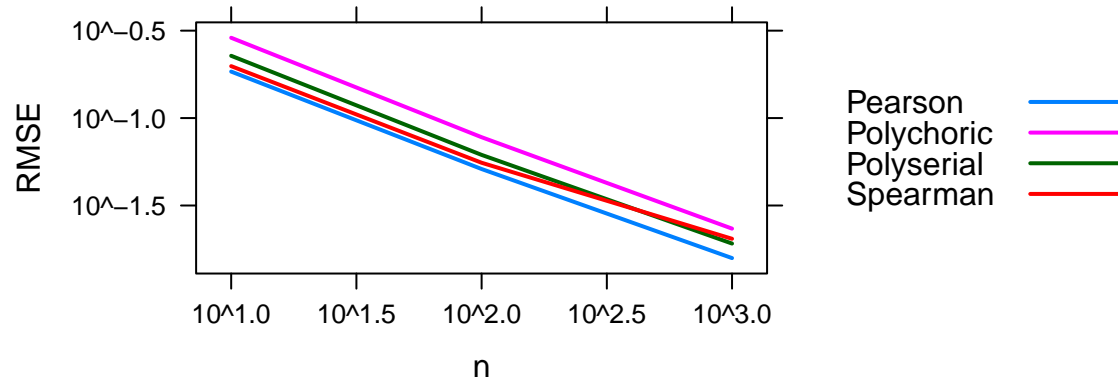
This plot shows the bias as a function of the true correlation  $\rho$ . Only the Spearman correlation shows a clear trend with a positive bias below 0 and a negative bias above zero.



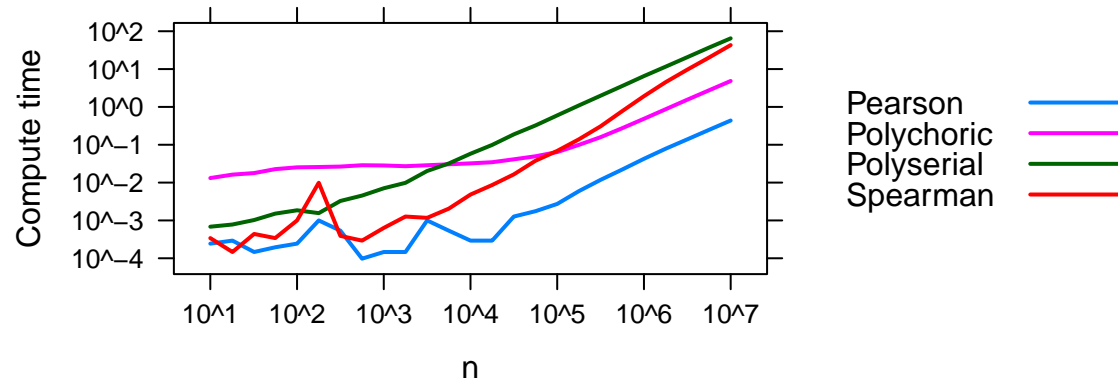
this plot shows the RMSE as a function of  $\rho$ . All of the correlation coefficients have a uniform RMSE as a function of  $\rho$  near  $\rho = 0$  that decreases near  $|\rho| = 1$ . All plots also show a decrease in RMSE as  $n$  increases. That is further shown in the next plot.



This plot shows the RMSE as a function of  $n$ . The plot shows a slope of about  $-\frac{1}{2}$ , which suggests first order consistency.



Finally, this plot shows the mean time to compute a single correlation coefficient as a function of  $\rho$  by  $n$  size.



## General setup for the weighted case

In this section the consistency of the weighted coefficients is explored.

For the weighted case the unweighted result would be consistent if there were not something about the higher weight cases that makes them different from the lower weight cases. Thus, while it is not reasonable to always assume that there is a difference between the high and low weight cases, both are possible and it serves as a more robust test of the methods in this package to consider a cases where they are associated. So, when sampling (other than simple random sampling with replacement) is used unweighted correlations may be consistent or may not be. The simulations below explore a case where the unweighted are not expected to be consistent to see if the weighted results are consistent.

## Weighted simulation study

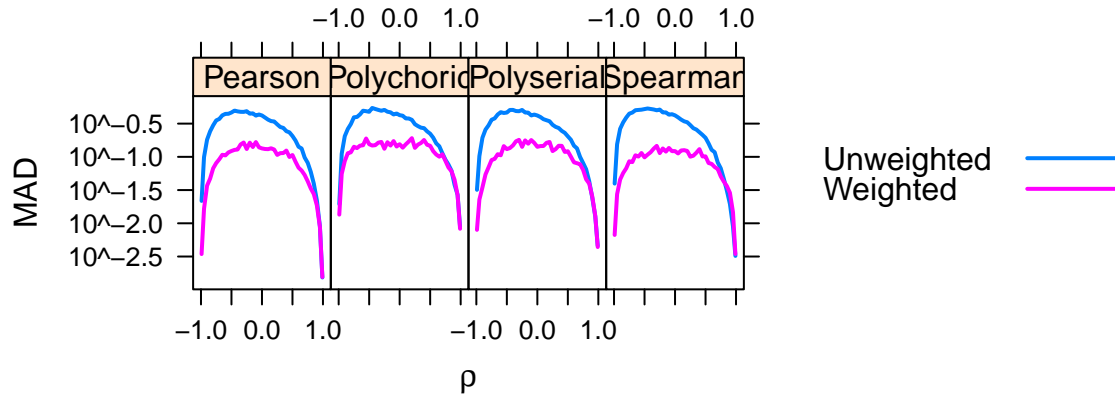
Simulations are carried out in the same fashion as before with a few extra steps to accommodate weights. The following changes were made:

- to find the correlation between  $n$  values,  $5n$  values were generated

- weights are assigned according to  $w_i = (x - y)^2 + 1$  and the probability of inclusion in the sample was proportional to  $Pr_i = \frac{\sum_j w_j}{w_i}$ .
- a sample of size  $n$  was then selected from the  $5n$  values according to the probabilities and weighted as well as unweighted correlation coefficients were calculated.

Two simulations were run. The first shows the MAD as a function of  $\rho$  and was run for  $n = 100$  and  $\rho \in (-0.99, -0.95, -0.90, -0.85, \dots, 0.95, 0.99)$ , with 100 iterations run for each value of  $\rho$ .

The following plot shows the MAD for the weighted and unweighted results as a function of  $\rho$  when  $n = 100$ .



The second simulation used the same values of  $\rho$  and used  $n \in \{10, 100, 1000, 10000, 100000\}$  and shows how RMSE and sample size are related. In particular, it shows first order convergence of the weighted correlation coefficient. Again, the fact that the simulations show that the unweighted correlation coefficient is not consistent is not meant to imply that it will always be that way—only that this is possible.

